

(12) INTERNATIONAL APPLICATION PUBLISHED UNDER THE PATENT COOPERATION TREATY (PCT)

(19) World Intellectual Property Organization
International Bureau



(43) International Publication Date
8 August 2002 (08.08.2002)

PCT

(10) International Publication Number
WO 02/061127 A2

(51) International Patent Classification⁷: C12Q 1/68

(74) Agent: GILL JENNINGS & EVERY; Broadgate House,
7 Eldon Street, London EC2M 7LH (GB).

(21) International Application Number: PCT/GB02/00439

(81) Designated States (national): AE, AG, AL, AM, AT, AU,

(22) International Filing Date: 30 January 2002 (30.01.2002)

AZ, BA, BB, BG, BR, BY, BZ, CA, CH, CN, CO, CR, CU,
CZ, DE, DK, DM, DZ, EC, EE, ES, FI, GB, GD, GE, GH,
GM, HR, HU, ID, IL, IN, IS, JP, KE, KG, KP, KR, KZ, LC,
LK, LR, LS, LT, LU, LV, MA, MD, MG, MK, MN, MW,
MX, MZ, NO, NZ, OM, PH, PL, PT, RO, RU, SD, SE, SG,
SI, SK, SL, TJ, TM, TN, TR, TT, TZ, UA, UG, US, UZ,
VN, YU, ZA, ZM, ZW.

(25) Filing Language: English

(84) Designated States (regional): ARIPO patent (GH, GM,

(26) Publication Language: English

KE, LS, MW, MZ, SD, SL, SZ, TZ, UG, ZM, ZW),
Eurasian patent (AM, AZ, BY, KG, KZ, MD, RU, TJ, TM),
European patent (AT, BE, CH, CY, DE, DK, ES, FI, FR,
GB, GR, IE, IT, LU, MC, NL, PT, SE, TR), OAPI patent
(BF, BJ, CF, CG, CI, CM, GA, GN, GQ, GW, ML, MR,
NE, SN, TD, TG).

(30) Priority Data:

09/771,708 30 January 2001 (30.01.2001) US

Published:

— without international search report and to be republished
upon receipt of that report

For two-letter codes and other abbreviations, refer to the "Guidance Notes on Codes and Abbreviations" appearing at the beginning of each regular issue of the PCT Gazette.

(72) Inventors; and

(75) Inventors/Applicants (for US only): BARNES, Colin [GB/GB]; Solexa Ltd., Chesterford Research Park, Little Chesterford, Nr. Saffron Walden, Essex CB10 1XL (GB). BALASUBRAMANIAN, Shankar [GB/GB]; University of Cambridge, Dept. of Chemistry, Lensfield Road, Cambridge CB2 1EW (GB). KLENERMAN, David [GB/GB]; University of Cambridge, Dept. of Chemistry, Lensfield Road, Cambridge CB2 1EW (GB).



WO 02/061127 A2

(54) Title: ARRAYED POLYNUCLEOTIDES AND THEIR USE IN GENOME ANALYSIS

(57) Abstract: A method for determining a single nucleotide polymorphism present in a genome, comprises the steps of: (i) immobilising fragments of said genome onto the surface of a solid support to form an array of polynucleotide molecules capable of interrogation, wherein the array allows the molecules to be individually resolved by optical microscopy, and wherein each molecule is immobilised by covalent bonding to the surface, other than at part of each molecule that can be interrogated; (ii) identifying nucleotides at selected positions in the genome; and (iii) comparing the results of step (ii) with a known consensus sequence, and identifying any differences between the consensus sequence and said genome.

ARRAYED POLYNUCLEOTIDES AND THEIR USE IN GENOME ANALYSIS

Field of the Invention

This invention relates to fabricated arrays of polynucleotides, and to their analytical applications. In particular, this invention relates to the use of fabricated polynucleotide arrays in methods for obtaining genetic sequence information.

Background of the Invention

Advances in the study of molecules have been led, in part, by improvement in technologies used to characterise the molecules or their biological reactions. In particular, the study of nucleic acids, DNA and RNA, has benefited from developing technologies used for sequence analysis and the study of hybridisation events.

An example of the technologies that have improved the study of nucleic acids, is the development of fabricated arrays of immobilised nucleic acids. These arrays typically consist of a high-density matrix of polynucleotides immobilised onto a solid support material. Fodor *et al.*, Trends in Biotechnology (1994) 12:19-26, describes ways of assembling the nucleic acid arrays using a chemically sensitised glass surface protected by a mask, but exposed at defined areas to allow attachment of suitably modified nucleotides. Typically, these arrays may be described as "many molecule" arrays, as distinct regions are formed on the solid support comprising a high density of one specific type of polynucleotide.

An alternative approach is described by Schena *et al.*, Science (1995) 270:467-470, where samples of DNA are positioned at predetermined sites on a glass microscope slide by robotic micropipetting techniques. The DNA is attached to the glass surface along its entire length by non-covalent electrostatic interactions. However, although hybridisation with complementary DNA sequences can occur, this approach may not permit the DNA to be freely available for interacting with other components such as polymerase enzymes, DNA-binding proteins etc.

Recently, the Human Genome Project determined the entire sequence of the human genome- all 3×10^9 bases. The sequence information represents that of an average human. However, there is still considerable interest in identifying differences in the genetic sequence between different individuals. The most common form of genetic variation is single nucleotide polymorphisms (SNPs). On average one base in 1000 is a SNP, which

means that there are 3 million SNPs for any individual. Some of the SNPs are in coding regions and produce proteins with different binding affinities or properties. Some are in regulatory regions and result in a different response to changes in levels of metabolites or messengers. SNPs are also found in non-coding regions, and these are also important as 5 they may correlate with SNPs in coding or regulatory regions. The key problem is to develop a low cost way of determining one or more of the SNPs for an individual.

The nucleic acid arrays may be used to determine SNPs, and they have been used to study hybridisation events (Mirzabekov, Trends in Biotechnology (1994) 12:27-32). Many of these hybridisation events are detected using fluorescent labels attached to 10 nucleotides, the labels being detected using a sensitive fluorescent detector, e.g. a charge-coupled detector (CCD). The major disadvantages of these methods are that it is not possible to sequence long stretches of DNA, and that repeat sequences can lead to ambiguity in the results. These problems are recognised in Automation Technologies for Genome Characterisation, Wiley-Interscience (1997), ed. T. J. Beugelsdijk, Chapter 10: 15 205-225.

In addition, the use of high-density arrays in a multi-step analysis procedure can lead to problems with phasing. Phasing problems result from a loss in the synchronisation of a reaction step occurring on different molecules of the array. If some of the arrayed molecules fail to undergo a step in the procedure, subsequent results obtained for these 20 molecules will no longer be in step with results obtained for the other arrayed molecules. The proportion of molecules out of phase will increase through successive steps and consequently the results detected will become ambiguous. This problem is recognised in the sequencing procedure described in US-A-5302509. This method is therefore not suitable for the determination of SNPs, where the precise identification of a particular 25 sequence is required.

WO-A-96/27025 is a general disclosure of single molecule arrays. Although sequencing procedures are disclosed, there is little description of the applications to which the arrays can be applied. There is also only a general discussion on how to prepare the arrays.

30 Summary of the Invention

According to one aspect of the present invention, a method for determining a single nucleotide polymorphism present in a genome, comprises the steps of:

5 (i) immobilising fragments of said genome onto the surface of a solid support to form an array of polynucleotide molecules capable of interrogation, wherein the array allows the molecules to be individually resolved by optical microscopy, and wherein each molecule is immobilised by covalent bonding to the surface, other than at that part of each molecule that can be interrogated;

10 (ii) identifying nucleotides at selected positions in the genome; and (iii) comparing the results of step (ii) with a known consensus sequence, and identifying any differences between the consensus sequence and said genome.

The arrays of the present invention comprise what are effectively single molecules. This has many important benefits for the study of the molecules and their interaction with other biological molecules. In particular, fluorescent labels can be used in interactions with the single polynucleotide molecules and can be detected using an optical microscope linked to a sensitive detector, resulting in a distinct signal for each polynucleotide.

The arrays permit a massively parallel approach to monitoring fluorescent or other events on the polynucleotides. Such massively parallel data acquisition makes the arrays extremely useful in the detection and characterisation of single nucleotide polymorphisms.

Description of the Invention

20 According to the present invention, the single polynucleotides immobilised onto the surface of a solid support should be capable of being resolved by optical means. This means that, within the resolvable area of the particular imaging device used, there must be one or more distinct images each representing one polynucleotide. Typically, the polynucleotides of the array are resolved using a single molecule fluorescence microscope 25 equipped with a sensitive detector, e.g. a charge-coupled device (CCD). Each polynucleotide of the array may be analysed simultaneously or, by scanning the array, a fast sequential analysis can be performed.

The polynucleotides of the array are derived from fragments of genomic DNA.

30 The term "single molecule" or "single polynucleotide" is used herein to distinguish from high density multi-molecule arrays in the prior art, which may comprise distinct clusters of many molecules of the same type.

The term "individually resolved" is used herein to indicate that, when visualised, it is possible to distinguish one polynucleotide on the array from its neighbouring polynucleotides. Visualisation may be effected by the use of reporter labels, e.g. fluorophores, the signal of which is individually resolved.

5 The terms "arrayed polynucleotides" and "polynucleotide arrays" are used herein to define a plurality of single polynucleotides. The term is intended to include the attachment of other molecules to a solid surface, the molecules having a polynucleotide attached that can be further interrogated during the SNP analysis. For example, the arrays 10 may comprise linker molecules immobilised on a solid surface, the linker molecules being conjugated or otherwise bound to a polynucleotide that may be interrogated, to determine the presence of a SNP.

The density of the array is not critical. However, the present invention can make use of a high density of single molecules (polynucleotides), and these are preferable. For example, arrays with a density of 10^6 - 10^9 polynucleotides per cm^2 may be used. 15 Preferably, the density is at least $10^7/\text{cm}^2$ and typically up to $10^8/\text{cm}^2$. These high density arrays are in contrast to other arrays which may be described in the art as "high density" but which are not necessarily as high and/or which do not allow single molecule resolution.

Using the methods and apparatus of the present invention, it may be possible to image at least 10^7 or 10^8 polynucleotides. Fast sequential imaging may be achieved using 20 a scanning apparatus; shifting and transfer between images may allow higher numbers of molecules to be imaged.

The extent of separation between the individual polynucleotides on the array will be determined, in part, by the particular technique used to resolve the individual polynucleotide. Apparatus used to image molecular arrays are known to those skilled in 25 the art. For example, a confocal scanning microscope may be used to scan the surface of the array with a laser to image directly a fluorophore incorporated on the individual molecule by fluorescence. Alternatively, a sensitive 2-D detector, such as a charge-coupled device, can be used to provide a 2-D image representing the individual polynucleotides on the array.

30 Resolving single polynucleotides on the array with a 2-D detector can be done if, at 100 x magnification, adjacent polynucleotides are separated by a distance of approximately at least 250nm, preferably at least 300nm and more preferably at least

350nm. It will be appreciated that these distances are dependent on magnification, and that other values can be determined accordingly, by one of ordinary skill in the art.

Other techniques such as scanning near-field optical microscopy (SNOM) are available which are capable of greater optical resolution, thereby permitting more dense 5 arrays to be used. For example, using SNOM, adjacent polynucleotides may be separated by a distance of less than 100nm, e.g. 10nm. For a description of scanning near-field optical microscopy, see Moyer *et al.*, *Laser Focus World* (1993) 29(10).

An additional technique that may be used is surface-specific total internal reflection 10 fluorescence microscopy (TIRFM); see, for example, Vale *et al.*, *Nature*, (1996) 380: 451-453). Using this technique, it is possible to achieve wide-field imaging (up to 100 $\mu\text{m} \times$ 100 μm) with single molecule sensitivity. This may allow arrays of greater than 10^7 resolvable polynucleotides per cm^2 to be used.

Additionally, the techniques of scanning tunnelling microscopy (Binnig *et al.*, *Helvetica Physica Acta* (1982) 55:726-735) and atomic force microscopy (Hansma *et al.*, 15 *Ann. Rev. Biophys. Biomol. Struct.* (1994) 23:115-139) are suitable for imaging the arrays of the present invention. Other devices which do not rely on microscopy may also be used, provided that they are capable of imaging within discrete areas on a solid support.

Single polynucleotides may be arrayed by immobilisation to the surface of a solid support. This may be carried out by any known technique, provided that suitable 20 conditions are used to ensure adequate separation. Generally the array is produced by dispensing small volumes of a sample containing a mixture of the fragmented genomic DNA onto a suitably prepared solid surface, or by applying a dilute solution to the solid surface to generate a random array. The formation of the array then permits interrogation of each arrayed polynucleotide to be carried out.

25 Suitable solid supports are available commercially, and will be apparent to the skilled person. The supports may be manufactured from materials such as glass, ceramics, silica and silicon. The supports usually comprise a flat (planar) surface, or at least an array in which the polynucleotides to be interrogated are in the same plane. Any suitable size may be used. For example, the supports might be of the order of 1-10 cm in each 30 direction.

Immobilisation may be by specific covalent or non-covalent interactions. Covalent attachment is preferred. Immobilisation will preferably be at either the 5' or 3' position,

so that the polynucleotide is attached to the solid support at one end only. However, the polynucleotide may be attached to the solid support at any position along its length, the attachment acting to tether the polynucleotide to the solid support. The immobilised polynucleotide is then able to undergo interactions at positions distant from the solid support. Typically the interaction will be such that it is possible to remove any molecules bound to the solid support through non-specific interactions, e.g. by washing. Immobilisation in this manner results in well separated single polynucleotides.

In one embodiment, the array comprises polynucleotides with a hairpin loop structure, one end of which comprises the target polynucleotide derived from the genomic DNA sample.

The term "hairpin loop structure" refers to a molecular stem and loop structure formed from the hybridisation of complementary polynucleotides that are covalently linked. The stem comprises the hybridised polynucleotides and the loop is the region that covalently links the two complementary polynucleotides. Anything from a 5 to 25 (or more) base pair double-stranded (duplex) region may be used to form the stem. In one embodiment, the structure may be formed from a single-stranded polynucleotide having complementary regions. The loop in this embodiment may be anything from 2 or more non-hybridised nucleotides. In a second embodiment, the structure is formed from two separate polynucleotides with complementary regions, the two polynucleotides being linked (and the loop being at least partially formed) by a linker moiety. The linker moiety forms a covalent attachment between the ends of the two polynucleotides. Linker moieties suitable for use in this embodiment will be apparent to the skilled person. For example, the linker moiety may be polyethylenè glycol (PEG).

There are many different ways of forming the hairpin structure to incorporate the target polynucleotide. However, a preferred method is to form a first molecule capable of forming a hairpin structure, and ligate the target polynucleotide to this. Ligation may be carried out either prior to or after immobilisation to the solid support. The resulting structure comprises the single-stranded target polynucleotide at one end of the hairpin and a primer polynucleotide at the other end.

The genomic DNA may be PCR-amplified or used directly to generate fragments of DNA using either restriction endonucleases, other suitable enzymes, a mechanical form of fragmentation or a non-enzymatic chemical fragmentation method. The fragments may

be of any suitable length, preferably from 20 to 2000 bases, more preferably 20 to 1000 bases, most preferably 20 to 200 bases. In the case of fragments generated by restriction endonucleases, hairpin structures bearing a complementary restriction site at the end of the first hairpin may be used, and selective ligation of one strand of the DNA sample 5 fragments may be achieved by various methods.

Method 1: The fragments are ligated to a hairpin made, for example, with a 3' overhang containing all possible sequences of a few nucleotides (preferably 3-20 bases long, more preferably 5-9 bases long), a 3' hydroxyl and a 5' phosphate. Ligation creates a 5' overhang that is capable of being sequenced from the 3' hydroxyl of the hairpin using 10 the newly ligated genomic fragment as a template by the methods described.

Method 2: in the design of the hairpin, a single (or more) base gap can be incorporated at the 3' end (the receded strand) such that upon ligation of the DNA fragment only one strand is covalently joined to the hairpin. The base gap can be formed by hybridising a further separate polynucleotide to the 5'-end of the first hairpin structure. 15 On ligation, the DNA fragment has one strand joined to the 5'-end of the first hairpin, and the other strand joined to the 3'-end of the further polynucleotide. The further polynucleotide (and the other strand of the DNA fragment) may then be removed by disrupting hybridisation.

Method 3: Genomic fragments are left in their double stranded-form or are made 20 to be double stranded and blunt ended by conventional means and are phosphatased to produce 3' and 5' hydroxyls as is known in the art. The fragments are ligated to a hairpin made for example with a blunt end, a 3' hydroxy and a 5' phosphate. Ligation of only one strand followed by denaturation and washing away of the unligated strand creates a 5' overhang that is capable of being sequenced from the 3' hydroxyl of the hairpin using the 25 newly ligated genomic fragment as a template by the methods described.

The net result should be covalent ligation of only one strand of a DNA fragment of genomic DNA, to the hairpin, the DNA fragment being then in the form of a 5' overhang that is capable of being sequenced. Such ligation reactions may be carried out in solution at optimised concentrations based on conventional ligation chemistry, for 30 example, carried out by DNA ligases or non-enzymatic chemical ligation. Should the fragmented DNA be generated by random shearing of genomic DNA, then the ends can be filled in with any polymerase to generate blunt-ended fragments which may be blunt-

end-ligated onto blunt-ended hairpins. Alternatively, the blunt-ended DNA fragments may be ligated to oligonucleotide adapters which are designed to allow compatible ligation with the sticky-end hairpins, in the manner described previously.

5 The hairpin-ligated DNA constructs may then be covalently attached to the surface of a solid support to generate the single molecule array, or ligation may follow attachment to form the array.

10 The arrays may then be used in procedures to determine the presence of a SNP. If the target fragments are generated via restriction digest of genomic DNA, the 15 recognition sequence of the restriction or other nuclease enzyme will provide 4, 6, 8 bases or more of known sequence (dependent on the enzyme). Further sequencing of at least 20 4 bases and preferably between 10 and 30 bases on the array should provide sufficient overall sequence information to place that stretch of DNA into unique context with a total human genome sequence, thus enabling the sequence information to be used for genotyping and more specifically single nucleotide polymorphism (SNP) scoring.

15 Simple calculations have suggested the following based on sequencing a 10^7 molecule array prepared from hairpin ligation: for a 6 base pair recognition sequence, a single restriction enzyme will generate approximately 10^6 ends of DNA. If a stretch of 13 bases is sequenced on the array (i.e. 13×10^6 bases), approximately 13,000 SNPs will be detected. The approach is therefore suitable for forensic analysis or any other system 20 which requires unambiguous identification of individuals to a level as low 10^3 SNPs.

It is of course possible to sequence the complete target polynucleotide, if required.

25 Sequencing can be carried out by the stepwise identification of suitably labelled nucleotides, referred to in US-A-5634413 as "single base" sequencing methods. The target polynucleotide is primed with a suitable primer (or prepared as a hairpin construct which will contain the primer as part of the hairpin), and the nascent chain is extended in a stepwise manner by the polymerase reaction. Each of the different nucleotides (A, T, G and C) incorporates a unique fluorophore which may be located at the 3' position to act 30 as a blocking group to prevent uncontrolled polymerisation. The polymerase enzyme incorporates a nucleotide into the nascent chain complementary to the target, and the blocking group prevents further incorporation of nucleotides. The array surface is then cleared of unincorporated nucleotides and each incorporated nucleotide is "read" optically by a charge-coupled detector using laser excitation and filters. The 3' -blocking group is

then removed (deprotected), to expose the nascent chain for further nucleotide incorporation.

Because the array consists of distinct optically resolvable polynucleotides, each target polynucleotide will generate a series of distinct signals as the fluorescent events are 5 detected. Details of the sequence are then determined and can be compared with known sequence information to identify SNPs.

The number of cycles that can be achieved is governed principally by the yield of the deprotection cycle. If deprotection fails in one cycle, it is possible that later deprotection and continued incorporation of nucleotides can be detected during the next 10 cycle. Because the sequencing is performed at the single molecule level, the sequencing can be carried out on different polynucleotide sequences at one time without the necessity for separation of the different sample fragments prior to sequencing. This sequencing also avoids the phasing problems associated with prior art methods.

The labelled nucleotides may comprise a separate label and removable blocking 15 group, as will be appreciated by those skilled in the art. In this context, it will usually be necessary to remove both the blocking group and the label prior to further incorporation.

Deprotection may be carried out by chemical, photochemical or enzymatic reactions. A similar, and equally applicable, sequencing method is disclosed in EP-A-0640146. Other suitable sequencing procedures will be apparent to the skilled person.

20 It is not necessary to determine the sequence of the full polynucleotide fragment. For example, it may be preferable to determine the sequence of 16-30 specific bases, which is sufficient to identify the DNA fragment by comparison to a consensus sequence, e.g. to that known from the Human Genome Project. Any SNP occurring within the sequenced region can then be identified. The specific bases do not have to be contiguous. For 25 example, the procedure may be carried out by the incorporation of non-labelled bases followed, at pre-determined positions, by the incorporation of a labelled base. Provided that the sequence of sufficient bases is determined, it should be possible to identify the fragment. Again, any SNPs occurring at the determined base positions, can be identified. For example, the method may be used to identify SNPs that occur after cytosine. 30 Template DNA (genomic fragments) can be contacted with each of the bases A, T and G, added sequentially or together, so that the complementary strand is extended up to a position that requires C. Non-incorporated bases can then be removed from the array,

followed by the addition of C. The addition of C is followed by monitoring the next base incorporation (using a labelled base). By repeating this process a sufficient number of times, a partial sequence is generated where each base immediately following a C is known. It will then be possible to identify the full sequence, by comparison of the partial sequence to a reference sequence. It will then also be possible to determine whether there 5 are any SNPs occurring after any C.

To further illustrate this, a device may comprise 10^7 restriction fragments per cm^2 . If 30 bases are determined for each fragment, this means 3×10^8 bases are identified. Statistically, this should determine 3×10^5 SNPs for the experiment. If the fragments each 10 comprise 1000 nucleotides, it is possible to have 10^{10} nucleotides per cm^2 , or three copies of the human genome. The approach therefore permits large sequence or SNP analysis to be performed.

The images and other information about the arrays, e.g. positional information, etc. are processed by a computer program which may perform image processing to reduce 15 noise and increase signal or contrast, as is known in the art. The computer program may perform an optional alignment between images and/or cycles, extract the single molecule data from the images, correlate the data between images and cycles and specify the DNA sequence from the patterns of signal produced from the individual molecules.

The individual DNA sequence reads of at least 4 bases, and more preferably at 20 least 16 bases in the case of human genomic DNA, and more preferably 16-30 bases, are aligned and compared with a genomic sequence. The methods for performing this alignment are based upon techniques known to those skilled in the art. The individual DNA sequence reads are aligned with respect to the reference sequence by finding the best match between the individual DNA sequence reads and the reference sequence. Using the 25 known alignments, one or many individual DNA sequence reads covering a given region of the genomic DNA sequence are obtained. All the aligned individual DNA sequence reads are interpreted at each nucleotide position in the reference sequence as either containing the identical sequence to the reference sequence, or containing an error in some of the individual DNA sequence reads, or containing a known or novel mutation, SNP, 30 deletion, insertion, etc. at that position. Furthermore, for most chromosomes, at each position in the reference sequence, the individual may contain one (homozygous) or two (heterozygous) different nucleotides corresponding to the two copies of each

chromosome. The sum total of all the individual variations in the reference sequence corresponding to a given individual sample is collectively referred to as a "total genotype".

The following Example illustrates the invention.

Example

5 Preparation of hairpin single molecule array (unlabelled DNA): A 10 μ M solution of oligonucleotide (5'-TCgACTgCTgAAAAgCgTCggCTggT-HEG-aminodT-HEG-ACCAgCCgACGCTTT; SEQ ID NO. 1) in DMF containing 10% water and 1% diisopropylethylamine (DIPEA) was prepared. To this, a stock solution of the GMBS crosslinker was added to give a final concentration of 1 mM N-[γ -10 Maleimidobutyryloxy]succinimide ester (GMBS) (100 eqvs.). The reaction was left for 1 h at room temperature, purified using a NAP size exclusion column and freeze-dried in aliquots that were re-dissolved immediately prior to use.

15 A fused silica slide was treated with decon for 12 h then rinsed with water, EtOH, dried and placed in a flow cell. A solution of the GMBS DNA (150 nM) and mercaptopropyltrimethoxysilane (3 μ M) in 9:1 sodium acetate (30 mM, pH 4.3): isopropanol was placed over the slide for 30 min. at 65 °C. The cell was flushed first with 50 mM Tris.HCl, 1 mM EDTA, pH 7.4 and then 50 mM Tris.HCl, 1 mM EDTA, 5 mM MgCl₂, 10 mM NaCl (pH 7.4) (10 mL) at 37 °C (TKF buffer). The cell was filled with 100 μ L of 2 μ M Cy5-dCTP, 2 μ M dTTP, 2 μ M dATP, 1 mM DTT, Klenow exo- (10 units) 20 in TKF buffer and incubated at 37 °C for 10 mins. then flushed with TKF buffer (20 mL) and TKF buffer containing NaCl (1 M) which removes bound protein. A second cycle consisting of 100 μ L of 2 μ M Cy3-dCTP, 2 μ M dGTP, 2 μ M dATP, 1 mM DTT, Klenow exo- (10 units) in TKF buffer was incubated at 37 °C for 10 mins. then flushed with TKF buffer (20 mL) and TKF buffer containing NaCl (1 M).

25 The flowcell was inverted so that the chamber coverslip contacts the objective lens of an inverted microscope (Nikon TE200) via an immersion oil interface. A 60° fused silica dispersion prism was optically coupled to the back of the slide through a thin film of glycerol. Laser light was directed at the prism such that at the glass/sample interface subtended an angle of approximately 68° to the normal of the slide and subsequently 30 underwent Total Internal Reflection (TIR). Fluorescence from the surface produced by excitation with the surface specific evanescent wave generated by TIR was collected by

the 100X objective lens of the microscope and imaged onto an intensified charged coupled device (ICCD) camera (Pentamax, Princeton Instruments).

Images were recorded using a combination of a 532 Nd:YAG laser with a 580DF30 emission filter (Omega optics) and a pumped dye laser at 630 nm with a 5 670DF40 emission filter. Images were recorded with an exposure of 500 ms and maximum camera gain and a laser power of 50 mW (green) and 40 mW (red) at the prism.

Two colour fluorophore labelled nucleotide incorporations were identified by the co-localisation of discreet points of fluorescence from single molecules of Cy3 and Cy5 following superimposing the two images. Molecules were considered co-localised when 10 fluorescent points were within a pixel separation of each other. For a 90 μ m and 90 μ m field projected onto a CCD array of 512'x 512 pixels the pixel size dimension is 176 nm.

An average 46.2% of Cy3 and 57.5% of Cy5 were colocalised; showing >50% of the molecules that underwent the Cy5 incorporation underwent a second cycle of Cy3 incorporation. In the absence of enzyme in the second cycle the level of Cy3 was greatly 15 reduced and the colocalisation was <2%. Polymerase fidelity controls, whereby the dATP or dGTP was omitted from the cycles, gave colocalisation levels of approximately 4%.

This demonstrates that sequence determination at the single molecule level can be achieved and makes it possible to extend this to genomic fragments to identify SNPs.

CLAIMS

1. A method for determining a single nucleotide polymorphism present in a genome, comprising
 - (i) immobilising fragments of said genome onto the surface of a solid support to form an array of polynucleotide molecules capable of interrogation, wherein the array allows the molecules to be individually resolved by optical microscopy, and wherein each molecule is immobilised by covalent bonding to the surface, other than at that part of each molecule that can be interrogated;
 - (ii) identifying nucleotides at selected positions in the genome; and
 - (iii) comparing the results of step (ii) with a known consensus sequence, and identifying any differences between the consensus sequence and said genome.
2. A method according to claim 1, wherein step (ii) comprises
 - (a) contacting the array with each of the bases A, T, G and C, under conditions that permit the polymerase reaction to proceed and thereby form sequences complementary to those in the array;
 - (b) determining the incorporation of a base at each of selected positions in the complementary sequences, optionally repeating steps (a) and (b).
3. A method according to claim 2, wherein each base contains a removable fluorescent label and each sequential base incorporation is determined.
4. A method according to claim 2 or claim 3, wherein each base contains a removable blocking group that prevents further base incorporation, and wherein the blocking group is removed after determining base incorporation.
5. A method according to claim 2, wherein (a) is carried out by first contacting the array with three of the bases under conditions that permit the polymerase reaction to proceed, removing unreacted bases from the array and incorporating the remaining-base, so that (b) proceeds only after incorporation of the remaining base.
6. A method according to any of claims 1 to 5, wherein adjacent polynucleotides of the array are separated by a distance of at least 10nm.
7. A method according to claim 6, wherein the polynucleotides are separated by a distance of at least 100nm.

8. A method according to claim 6, wherein the polynucleotides are separated by a distance of at least 250nm.
9. A method according to any preceding claim, wherein the array has a density of from 10^6 to 10^9 polynucleotides per cm^2 .
- 5 10. A method according to claim 9, wherein the density is from 10^7 to 10^8 polynucleotides per cm^2 .
11. A method according to any preceding claim, wherein the polynucleotides are immobilised to the solid support via the 5' terminus, the 3' terminus or via an internal nucleotide.

SEQUENCE LISTING

<110> Solexa Ltd.

<120> Arrayed Polynucleotides and Their Use in Genome
Analysis

<130> REP07019WO

<140> not yet known

<141> 2002-01-30

<150> 09/771708

<151> 2001-01-30

<160> 1

<170> PatentIn Ver. 2.1

<210> 1

<211> 42

<212> DNA

<213> Artificial Sequence

<220>

<221> misc_feature

<222> (1)..(42)

<223> n = hexaethyleneglycol-aminodT-hexaethyleneglycol

<220>

<223> Description of Artificial Sequence: Synthetic
oligonucleotide

<400> 1

tcgactgctg aaaagcgtcg gctggtnacc agccgacgct tt

42